

1 A METHOD FOR ARBITRATING BANDWIDTH IN A COMMUNICATIONS SWITCH

2
3 BACKGROUND OF THE INVENTION
4

5 1. Field of the Invention

6 The invention relates to telecommunications networks. More
7 particularly, the invention relates to a method for arbitrating
8 bandwidth in a telecommunications network switch.

9
10 2. State of the Art

11 One of the earliest techniques for employing broadband
12 telecommunications networks was called time division multiplexing
13 (TDM). The basic operation of TDM is simple to understand. A
14 high frequency signal is divided into multiple time slots within
15 which multiple lower frequency signals can be carried from one
16 point to another. The actual implementation of TDM is quite
17 complex, however, requiring sophisticated framing techniques and
18 buffers in order to accurately multiplex and demultiplex signals.
19 The North American standard for TDM (known as T1 or DS1) utilizes
20 twenty-four interleaved channels together having a rate of 1.544
21 Mbits/sec. The European standard for TDM is known as E-1 and
22 utilizes thirty interleaved channels having a rate of 2.048
23 Mbits/sec. A hierarchy of multiplexing is based on multiples of
24 the T1 or E-1 signal, one of the most common being T3 or DS3. A
25 T3 signal has 672 channels, the equivalent of twenty-eight T1

1 signals. TDM was originally designed for voice channels. Today,
2 however, it is used for both voice and data.

3
4 An early approach to broadband data communication was called
5 packet switching. One of the differences between packet switching
6 and TDM is that packet switching includes methods for error
7 correction and retransmission of packets which become lost or
8 damaged in transit. Another difference is that, unlike the
9 channels in TDM, packets are not necessarily fixed in length.
10 Further, packets are directed to their destination based on
11 addressing information contained within the packet. In contrast,
12 TDM channels are directed to their destination based on their
13 location in the fixed frame. Today, a widely used packet
14 switching protocol is known as IP (Internet Protocol).

15
16 More recently, broadband technologies known as ATM and SONET
17 have been developed. The ATM network is based on fixed length
18 packets (cells) of 53-bytes each (48-bytes payload with 5-bytes
19 overhead). One of the characteristics of the ATM network is that
20 users contract for a quality of service (QOS) level. Thus, ATM
21 cells are assigned different priorities based on QOS. For
22 example, constant bit rate (CBR) service is the highest priority
23 service and is substantially equivalent to a provisioned TDM
24 connection. Variable bit rate (VBR) service is an intermediate
25 priority service which permits the loss of cells during periods of

1 congestion. Unspecified bit rate (UBR) service is the lowest
2 priority and is used for data transmission which can tolerate high
3 latency such as e-mail transmissions.

4
5 The SONET network is based on a frame of 810-bytes within
6 which a 783-byte synchronous payload envelope (SPE) floats. The
7 payload envelope floats because of timing differences throughout
8 the network. The exact location of the payload is determined
9 through a relatively complex system of stuffs/destuffs and
10 pointers. In North America, the basic SONET signal is referred to
11 as STS-1 (or OC-1). The SONET network includes a hierarchy of
12 SONET signals wherein up to 768 STS-1 signals are multiplexed
13 together providing the capacity of 21,504 T1 signals (768 T3
14 signals). STS-1 signals have a frame rate of 51.84 Mbit/sec, with
15 8,000 frames per second, and 125 microseconds per frame. In
16 Europe, the base (STM-1) rate is 155.520 Mbit/sec, equivalent to
17 the North American STS-3 rate ($3 \times 51.84 = 155.520$), and the payload
18 portion is referred to as the virtual container (VC). To
19 facilitate the transport of lower-rate digital signals, the SONET
20 standard uses sub-STS payload mappings, referred to as Virtual
21 Tributary (VT) structures. (The ITU calls these Tributary Units or
22 TUs.) Four virtual tributary sizes are defined: VT-1.5, VT-2,
23 VT-3 and VT-6. VT-1.5 has a data transmission rate of 1.728
24 Mbit/s and accommodates a T1 signal with overhead. VT-2 has a
25 data transmission rate of 2.304 Mbit/s and accommodates an E1

1 signal with overhead. VT-3 has a data transmission rate of 3.456
2 Mbit/s and accommodates a T2 signal with overhead. VT-6 has a
3 data transmission rate of 6.912 Mbit/s and accommodates a DS2
4 signal with overhead.

5
6 Each of the above described broadband technologies can be
7 categorized as TDM, ATM, or Packet technologies, with SONET being
8 a complex form of TDM. From the foregoing, it will be
9 appreciated that TDM, ATM and Packet each have their own unique
10 transmission requirements. Consequently, different kinds of
11 switches are used to route these different kinds of signals. In
12 particular, TDM requires careful time synchronization; ATM
13 requires careful attention to the priority of cells and QOS; and
14 packet (e.g. IP) requires the ability to deal with variable length
15 packets. For these reasons, switching technologies for TDM, ATM,
16 and variable length packet switching have evolved in different
17 ways. Service providers and network designers have thus been
18 forced to deal with these technologies separately, often providing
19 overlapping networks with different sets of equipment which can
20 only be used within a single network.

SUMMARY OF THE INVENTION

It is therefore an object of the invention to provide methods and apparatus whereby different kinds of broadband signals can be switched through a single switching fabric.

It is also an object of the invention to provide a network element which can switch TDM, ATM, and variable length packet traffic all through the same switch fabric.

It is another object of the invention to provide a network switch chipset which can be combined with identical chip sets to provide a scalable network switch fabric.

It is a further object of the invention to provide a network switch which allows flexible partitioning among TDM, ATM, and variable length packet traffic.

Another object of the invention is to provide a network switch with redundant switch planes so that the failure of switch elements or links does not immediately cause a connection failure.

A further object of the invention is to provide a network switch which handles multicast as well as unicast voice and data transmission.

1
2 An additional object of the invention to provide a network
3 switch which supports Clos architectures as well as folded Clos
4 architectures.

5
6 In accord with these objects which will be discussed in
7 detail below, the network switch of the present invention includes
8 at least one port processor (also referred to in the appendices as
9 a "service processor") and at least one switch element. The port
10 processor has a SONET OC-x (SONET/SDH STS-x/STM-y) interface (for
11 TDM traffic), a UTOPIA and UTOPIA-frame based interface (for ATM
12 and packet traffic), and an interface to the switch element. An
13 exemplary port processor has a total I/O bandwidth equivalent to a
14 SONET OC-48 signal. An exemplary switch element has 12x12 ports
15 and supports a total bandwidth of 30 Gbps.

16
17 A typical switch according to the invention includes multiple
18 port processors and multiple switch elements. For a 48x48
19 "folded" switch, 48 port processors are coupled (four each) to 12
20 (first and third stage) switch elements and each of these twelve
21 switch elements is coupled to 8 (second stage) switch elements. A
22 three stage non-blocking switch according to the invention
23 provides a total bandwidth of 240 Gbps and a five stage non-
24 blocking switch provides a total bandwidth of 1 Tbps. An
25 exemplary three stage folded Clos architecture switch includes

1 forty-eight port processors and twenty switch elements. Four port
2 processors are coupled to each of twelve (first and third stage)
3 switch elements. Each of the twelve (first and third stage)
4 switch elements are coupled to eight (second stage) switch
5 elements. According to the presently preferred embodiment, each
6 port processor is provided with means for coupling to two ports of
7 a switch element or one port of two switch elements thereby
8 providing redundancy in the event of a link failure.

9
10 According to the invention, a data frame of 9 rows by 1700
11 slots is used to transport ATM, TDM, and Packet data from a port
12 processor through one or more switch elements to the same or
13 another port processor. Each frame is transmitted in 125
14 microseconds, each row in 13.89 microseconds. Each slot includes
15 a four-bit tag plus a four-byte payload (i.e., thirty-six bits).
16 The slot bandwidth (1/1700 of the total frame) is 2.592 Mbps which
17 is large enough to carry an E-1 signal with overhead. The four-
18 bit tag is a cross connect pointer which is set up when a TDM
19 connection is provisioned. The last twenty slots of the frame are
20 reserved for link overhead. Thus, the frame is capable of
21 carrying the equivalent of 1,680 E-1 TDM signals even though an
22 STM-16 frame has a capacity of only 1008 E-1 signals.

23
24 For ATM and packet data, a PDU (protocol data unit) of
25 sixteen slots is defined for a sixty-four-byte payload (large

1 enough to accommodate an ATM cell with switch overhead). A
2 maximum of ninety-six PDUs per row is permitted. The sixteen
3 four-bit tags of a PDU are not needed for PDU routing so they are
4 used as parity bits to protect the ATM or variable length packet
5 payload. Of the sixty-four-byte payload, twelve bytes (96 bits)
6 are used by the switch for internal routing. This leaves fifty-
7 two bytes for actual payload which is sufficient to carry an ATM
8 cell (without the one-byte HEC) and sufficient for larger packets
9 after fragmentation. The PDUs are self-routed through the switch
10 with a twenty-eight-bit routing tag which allows routing through
11 seven switch stages using four-bits per stage. The remaining
12 sixty-eight bits of the PDU are used for various other addressing
13 information such as indicating whether the PDU contains an ATM
14 cell, a packet, or a control message, whether reassembly of the
15 packet should be aborted, whether the payload is a first fragment,
16 middle fragment or last fragment, how many payload bytes are in
17 the last fragment, the fragment sequence count, and a destination
18 flow identifier.

19
20 The link overhead (LOH) in the last twenty slots of the frame
21 is analogous in function to the line and section overhead in a
22 SONET frame. The LOH may contain a 36-bit frame alignment pattern
23 which is used to delineate the byte and row boundaries from serial
24 data streams, a 32-bit status register for each output link, a 32-
25 bit switch and link identifier, and a 32-bit stuff pattern.

1 Since ATM and Packet traffic are typically not provisioned,
2 bandwidth must be arbitrated among ATM and Packet connections as
3 traffic enters the system. Moreover, since TDM traffic shares the
4 same frame as ATM and Packet traffic, bandwidth must be arbitrated
5 while maintaining TDM timing. According to the invention,
6 bandwidth is arbitrated by a system of requests and grants which
7 is implemented for each PDU in each row of the frame. The switch
8 elements provide three channels per link, two of which are used to
9 carry data and arbitration requests and one of which is used to
10 carry arbitration grants. According to the presently preferred
11 embodiment, a forty-eight-bit (1.5 slot) request element is
12 generated for each PDU in the next row of the frame. Each switch
13 element includes a single request parser and a separate request
14 arbitration module for each output link. The request elements are
15 generated by the port processors and include intra-switch "hop-by-
16 hop" routing tags and priority level information. Request
17 elements are buffered by the switch elements and low priority
18 request elements are discarded by a switch element if the buffer
19 fills. Each request element which is not discarded as it travels
20 through the switch fabric is returned to the port processor from
21 which it originated during one "row time", i.e. 13.89
22 microseconds. As suggested above, requests are made "in band"
23 interleaved with data and grants (the returned request elements)
24 are made "out of band" using the third channel in each link.

25

1 In order to maintain timing for TDM traffic, the V1-V4 bytes
2 in the VT/VC frame are stripped off and the VC bytes are buffered
3 at ingress to the switch by a port processor. The V1-V4 bytes are
4 regenerated at the egress from the switch by a port processor. In
5 rows having both PDU and TDM traffic, the PDUs are configured
6 early and the TDM slots are configured late in the row.

7
8 According to the presently preferred embodiment, each switch
9 element includes a multicast controller and a separate multicast
10 PDU buffer. Multicast request elements flow through the switch in
11 the same manner as standard unicast request elements. At the
12 point where the message needs to be multicast, the hop-by-hop
13 field's bit code for that switch stage indicates that the request
14 is multicast. The request is forwarded to the multicast
15 controller. On the grant path, the multicast controller sources a
16 grant if there is room for the data in the multicast
17 recirculating buffers. Once the data has been transmitted to the
18 multicast buffer, the multicast controller examines the data
19 header and determines on which output links it needs to be sent
20 out. At this point, the multicast controller sources a number of
21 request messages which are handled in the same manner as unicast
22 requests.

23
24 Additional objects and advantages of the invention will
25 become apparent to those skilled in the art upon reference to the

1 detailed description taken in conjunction with the provided
2 figures.

3
4 BRIEF DESCRIPTION OF THE DRAWINGS

5
6 Figure 1 is a simplified schematic diagram of a port
7 processor according to the invention;

8
9 Figure 2 is a simplified schematic diagram of a switch
10 element according to the invention;

11
12 Figure 3 is a schematic diagram illustrating the data frame
13 structure of the invention;

14
15 Figure 3a is a schematic diagram illustrating the presently
16 preferred format of a PDU according to the invention;

17
18 Figure 3b is a schematic diagram illustrating the row
19 structure including request elements to a first stage of the
20 switch;

21
22 Figure 3c is a schematic diagram illustrating the row
23 structure including request elements to a second stage of the
24 switch;

1 Figure 4 is a schematic illustration of a three stage 48x48
2 switch according to the invention; and

3
4 Figure 5 is a schematic illustration of a 48x48 folded Clos
5 architecture switch according to the invention.

6
7 BRIEF DESCRIPTION OF THE APPENDIX

8
9 Appendix A is an engineering specification (Revision 0.3) for
10 a port processor according to the invention; and

11
12 Appendix B is an engineering specification (Revision 0.3) for
13 a switch element according to the invention.

14
15 DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENTS

16
17 The apparatus of the invention generally includes a port
18 processor and a switch element. Figure 1 illustrates the main
19 features of the port processor 10, and Figure 2 illustrates the
20 main features of the switch element 100. Referring now to Figure
21 1, the port processor 10 includes a SONET interface and a UTOPIA
22 interface. On the ingress (RX) side, the SONET interface includes
23 a serial to parallel converter 12, a SONET framer and transport
24 overhead (TOH) extractor 14, a high order pointer processor 16,
25 and a path overhead (POH) extractor 18. For ATM and IP packets

1 transported in an SPE, the ingress side of the SONET interface
2 includes forty-eight HDLC framers 20 (for IP), forty-eight cell
3 delineators 22 (for ATM), and forty-eight 64-byte FIFOs 24 (for
4 both ATM and IP. For TDM signals transported in an SPE, the
5 ingress side of the SONET interface includes a demultiplexer and
6 low order pointer processor 26. On the egress (TX) side, the
7 SONET interface includes, for TDM signals, a multiplexer and low
8 order pointer generator 28. For ATM and IP packets transported
9 in an SPE, the egress side of the SONET interface includes forty-
10 eight 64-byte FIFOs 30, forty-eight HDLC frame generators 32, and
11 forty-eight cell mappers 34. The egress side of the SONET
12 interface also includes a POH generator 36, a high order pointer
13 generator 38, a SONET framer and TOH generator 40, and a parallel
14 to serial interface 42. On the ingress side, the UTOPIA interface
15 includes a UTOPIA input 44 for ATM and Packets and one 4x64-byte
16 FIFO 46. On the egress side, the UTOPIA interface includes
17 ninety-six 4x64-byte FIFOs 48 and a UTOPIA output 50.

18
19 The ingress portion of the port processor 10 also includes a
20 switch mapper 52, a parallel to serial switch fabric interface 54,
21 and a request arbitrator 56. The egress portion of the port
22 processor also includes a serial to parallel switch fabric
23 interface 58, a switch demapper 60, and a grant generator 62.

1 For processing ATM and packet traffic, the port processor 10
2 utilizes, at the ingress portion, a descriptor constructor 64, an
3 IPF and ATM lookup processor 66, an IP classification processor
4 68, an RED/Policing processor 70, all of which may be located off-
5 chip. These units process ATM cells and packets before handing
6 them to a (receive) data link manager 72. At the egress portion
7 of the port processor, a (transmit) data link manager 74 and a
8 transmit scheduler and shaper 76 are provided. Both of these
9 units may be located off-chip. The port processor is also
10 provided with a host interface 78 and a weighted round robin
11 scheduler 80.

12
13 The purpose of the port processor at ingress to the switch is
14 to unpack TDM, Packet, and ATM data and frame it according to the
15 data frame described below with respect to Figure 3. The port
16 processor also buffers TDM and packet data while making
17 arbitration requests for link bandwidth through the switch element
18 and grants arbitration requests received through the switch as
19 described in more detail below. In order to maintain timing for
20 TDM traffic, the V1-V4 bytes in the SONET frame are stripped off
21 and the VC bytes are buffered at the ingress to the switch. In
22 rows having both PDU and TDM traffic, it is preferable that the
23 PDUs are configured early and the TDM slots are configured late in
24 the row. At the egress of the switch, the port processor

1 reassembles TDM, Packet, and ATM data. The V1-V4 bytes are
2 regenerated at the egress from the switch.

3
4 Though not shown in Figure 1, the port processor 10 includes
5 dual switch element interfaces which permit it to be coupled to
6 two switch elements or to two ports of one switch element. When
7 both interfaces are used, the "standby" link carries only frame
8 information until a failure in the main link occurs and then data
9 is sent via the standby link. This provides for redundancy in the
10 switch so that connections are maintained even if a portion of the
11 switch fails.

12
13 Turning now to Figure 2, a switch element 100 according to
14 the invention includes twelve "datapath and link bandwidth
15 arbitration modules" 102 (shown only once in Figure 2 for
16 clarity). Each module 102 provides one link input 104 and one
17 link output 106 through the switch element 100. Those skilled in
18 the art will appreciate that data entering any link input can,
19 depending on routing information, exit through any link output.
20 According to the invention, each module 102 provides two forward
21 datapaths 108, 110, 112, 114 and one return "grant" path 116, 118.
22 The three paths are collectively referred to as constituting a
23 single channel. The reason why two datapaths are provided is to
24 increase the bandwidth of each channel. The two datapaths are
25 interleaved to provide a single "logical" serial datastream which

1 exceeds (doubles) the bandwidth of a single physical datastream.
2 Data is routed from an input link 104 to an output link 106 via an
3 input link bus 120 and an output link bus 122. Return path grants
4 are routed from an output link 106 to an input link 104 via a
5 grant bus 124.

6
7 The forward datapaths of each "datapath and link bandwidth
8 arbitration module" 102 include a data stream deserializer 126, a
9 data stream demapper 128, a row buffer mapper 130, a row buffer
10 132, a request arbitration module 134, a data stream mapper 136,
11 and a data stream serializer 138. The return grant path for each
12 module 102 includes a grant stream deserializer 140, a grant
13 stream demapper 142, a grant arbitration module 144, a grant
14 stream mapper 146, and a grant stream serializer 148.

15
16 The switch element 100 also includes the following modules
17 which are instantiated only once and which support the functions
18 of the twelve "datapath and link bandwidth arbitration modules"
19 102: a link synchronization and timing control 150, a request
20 parser 152, a grant parser 154, and a link RISC processor 156.
21 The switch element 100 also includes the following modules which
22 are instantiated only once and which support the other modules,
23 but which are not directly involved in "switching": a
24 configuration RISC processor 158, a system control module 160, a
25 test pattern generator and analyzer 162, a test interface bus

1 multiplexer 164, a unilink PLL 166, a core PLL 168, and a JTAG
2 interface 170.

3
4 A typical switch according to the invention includes multiple
5 port processors 10 and multiple switch elements 100. For
6 example, as shown in Figure 4, forty-eight "input" port
7 processors are coupled to twelve "first stage" switch elements,
8 four to each. Each of the first stage switch elements is coupled
9 to eight second stage switch elements. Each of the second stage
10 switch elements is coupled to twelve third stage switch elements.
11 Four "output" port processors are coupled to each of the third
12 stage switch elements. From the foregoing, those skilled in the
13 art will appreciate that the port processors and the switch
14 elements of invention can be arranged in a folded Clos
15 architecture as shown in Figure 5 where a single switch element
16 acts as both first stage and third stage.

17
18 Before describing in detail the functions of the port
19 processor 10 and the switch element 100, it should be appreciated
20 that the invention utilizes a unique framing technique which is
21 well adapted to carry combinations of TDM, ATM, and Packet data in
22 the same frame. Turning now to Figure 3, according to the
23 invention, a data frame of nine rows by 1700 slots is used to
24 transport ATM, TDM, and Packet data from a port processor through
25 one or more switch elements to a port processor. Each frame is

1 transmitted in 125 microseconds, each row in 13.89 microseconds.
2 Each slot includes a four-bit tag plus a four-byte payload (i.e.,
3 thirty-six bits). The slot bandwidth (1/1700 of the total frame)
4 is 2.592 Mbps which is large enough to carry an E-1 signal with
5 overhead. The four-bit tag is a cross connect pointer which is
6 set up when a TDM connection is provisioned. The last twenty
7 slots of the frame are reserved for link overhead (LOH). Thus,
8 the frame is capable of carrying the equivalent of 1,680 E-1 TDM
9 signals. The link overhead (LOH) in the last twenty slots of the
10 frame is analogous in function to the line and section overhead in
11 a SONET frame.

12
13 The contents of the LOH slots are inserted by the switch
14 mapper (52 in Figure 1). There are four types of data which may
15 be inserted in the LOH slots. A 36-bit framing pattern is
16 inserted into one of the twenty slots. The framing pattern is
17 common to all output links and configurable via a software
18 programmable register. A 32-bit status field is inserted into
19 another slot. The status field is unique for each output link and
20 is configurable via a software programmable register. A 32-bit
21 switch and link identifier is inserted into another slot. The
22 switch and link identifier includes a four bit link number, a
23 twenty-four bit switch element ID, and a four bit stage number. A
24 32-bit stuff pattern is inserted into slots not used by framing,

1 status, or ID. The stuff pattern is common to all output links
2 and is configurable via a software programmable register.

3
4 For ATM and packet data, a PDU (protocol data unit) of
5 sixteen slots is defined for a sixty-four-byte payload (large
6 enough to accommodate an ATM cell with overhead). The format of
7 the PDU is illustrated in Figure 3a. A maximum of ninety-six PDUs
8 per row is permitted (it being noted that the maximum number of
9 ATM cells in a SONET OC-48 row is seventy-five). The sixteen
10 four-bit tags (bit positions 32-35 in each slot) are not needed
11 for PDU routing so they are used as parity bits to protect the ATM
12 or IP payload. Of the sixty-four-byte payload, twelve bytes (96
13 bits) are used by the switch for internal routing (slots 0-2, bit
14 positions 0-31). This leaves fifty-two bytes (slots 3-15) for
15 actual payload which is sufficient to carry an ATM cell (without
16 the one-byte HEC) and sufficient for larger packets after
17 fragmentation. The PDUs are self-routed through the switch with a
18 twenty-eight-bit routing tag (slot 0, bit positions 0-27) which
19 allows routing through seven stages using four bits per stage.
20 The remaining sixty-eight bits of the PDU are used for various
21 other addressing information.

22
23 As shown in Figure 3a, the PDU bits at slot 0, bits 30-31 are
24 used to identify whether the PDU is idle (00), an ATM cell (01),
25 an IP packet (10), or a control message (11). The two bits at

1 slot 1, bit positions 30-31 are used to indicate the internal
2 protocol version of the chip which produced the PDU. For Packets
3 and control messages, the "valid bytes" field (slot 1, bits 24-29)
4 are used to indicate how many payload bytes are carried by the PDU
5 when the FragID field indicates that the PDU is the last fragment
6 of a fragmented packet. The VOQID field (slot 1, bit positions
7 19-23) identifies the class of service for the PDU. The class of
8 service can be a value from 0 to 31, where 0 is the highest
9 priority and 31 is the lowest. The FragID at slot 1, bits 17-18
10 indicates whether this PDU is a complete packet (11), a first
11 fragment (01), a middle fragment (00), or a last fragment (10).
12 The A bit at slot 1, bit position 16 is set if reassembly for this
13 packet is being aborted, e.g. because of early packet (or partial
14 packet) discard operations. When this bit is set, fragments of
15 the packet received until this point are discarded by the output
16 port processor. The fields labelled FFS are reserved for future
17 use. The Seq# field at slot 1, bits 0-3 is a modular counter
18 which counts packet fragments. The DestFlowId field at slot 2,
19 bits 0-16 identifies the "flow" in the destination port processor
20 to which this PDU belongs. A "flow" is an active data connection.
21 There are 128K flows per port processor.

22
23 As mentioned above, since ATM and Packet traffic are
24 typically not provisioned, bandwidth must be arbitrated among ATM
25 and Packet connections as traffic enters the system. Moreover,

1 since TDM traffic shares the same frame as ATM and Packet traffic,
2 bandwidth must be arbitrated while maintaining TDM timing.
3 According to the invention, bandwidth is arbitrated by a system of
4 requests and grants which is implemented for each PDU in each row
5 of the frame. The request elements, which are generated by the
6 port processors include "hop-by-hop" internal switch routing tags,
7 switch element stage, and priority information. According to the
8 presently preferred embodiment two request elements are sent in a
9 three contiguous slot bundle and at least eight slots of non-
10 request element traffic must be present between request element
11 bundles. The time separation between request element bundles is
12 used by the arbitration logic in the switch elements and the port
13 processors to process the request elements. The presently
14 preferred request element format is shown in section 7.1.5 of
15 Appendix B.

16
17 Figure 3b illustrates one example of how the row slots may be
18 allocated for carrying PDUs and request elements. As shown,
19 the maximum PDU capacity for a row is ninety-six. A block of
20 sixteen slots which is capable of carrying a single PDU is
21 referred to as a "group". For each group in the row, 1.5 slots of
22 bandwidth are required for carrying a forty-eight-bit request
23 element (RE). Figure 3b illustrates how two REs are inserted into
24 three slots within each of the first twenty-four groups. All the
25 REs should be carried within the row as early as possible in order

1 to allow the REs to ripple through the multistage switch fabric as
2 soon as possible after the start of a row. Section 7 of Appendix
3 B explains in detail how this affects the arbitration process.

4
5 The structure shown in Figure 3b is presently considered to
6 be the optimal format (for the first link) given system
7 requirements and implementation constraints. It places the REs
8 early in the row but spaces them out enough to allow for
9 arbitration. According to the presently preferred embodiment, the
10 row structure is somewhat different depending on for which link of
11 the switch it is configured. Figure 3b represents the row
12 structure between the port processor and a switch element of the
13 first switch fabric stage. The first block of two REs occupy the
14 first three slots of the row. The present implementation of the
15 arbitration logic which processes REs requires at least twelve
16 slot times of latency between each three-slot block of REs on the
17 input link. Also, there must be some latency from when the first
18 REs of the row are received by a switch element to when the REs
19 are inserted into the output link of the switch element. This
20 latency is used by the arbitration logic for mapping incoming REs
21 into the RE buffers. Thus, the row structure for the link between
22 the first stage and the second stage should have the first group
23 of REs starting at slot time 32. This is illustrated in Figure 3c
24 which shows the same structure as Figure 3b offset by thirty-two
25 slot times.

1
2 According to the presently preferred embodiment, TDM traffic
3 may be switched through the switch elements with a finest
4 granularity of one slot per row. The TDM traffic is switched
5 through the same path for a given slot for every row. The switch
6 elements will not allow different switch paths for the same TDM
7 data slot for different rows within the frame. This means that
8 the switch does not care about what the current row number is
9 (within a frame). The only time row numbering matters is when
10 interpreting the contents of the Link Overhead slots.

11
12 With a finest granularity of one slot per row, the switch
13 elements can switch TDM traffic with a minimum of 2.52 Mbps of
14 switching bandwidth. Since a slot can carry the equivalent of
15 four columns of traffic from a SONET SPE, it can be said that the
16 switch elements switch TDM traffic with a granularity of a VT1.5
17 or VT2 channel. Although a VT1.5 channel only occupies three
18 columns in the SONET SPE, it will still be mapped to the slot
19 format which is capable of holding four SPE columns. As mentioned
20 above, the format of the contents of the thirty-six-bit slot
21 carrying TDM traffic is a four-bit tag and a thirty-two bits of
22 payload. The tag field definitions are shown in Table 1 below.

0000	Idle
0001	reserved
1010	reserved
1011	Data present
1100	V5 byte in bits 31-24
1101	V5 byte in bits 23-16
1110	V5 byte in bits 15-8
1111	V5 byte in bits 7-0

Table 1

The switch elements know whether or not a slot contains TDM data via preconfigured connection tables. These tables are implemented as an Input Cross Connect RAM for each input link. The input slot number is the address into the RAM, while the data output of the RAM contains the destination output link and slot number. The connection table can be changed by a centralized system controller which can send control messages to the switch elements via either of two paths: (1) a host interface port or (2) in-band control messages which are sent via the link data channel. Since TDM connections will be changed infrequently, this relatively slow control message approach to update the connection tables is acceptable. It is the responsibility of an external software module to determine and configure the connection tables within the switch elements such that no TDM data will be lost.

1 Returning now to Figure 1, the receive side SONET interface
2 of the port processor 10 includes the deserializer 12 and framer
3 14. This interface may be configured as one OC-48, 16-bits wide
4 at 155 MHz, four OC-12s, serially at 622 MHz, or four OC-3s,
5 serially at 155 MHz. When configured as one OC-48, the
6 deserializer 12 is not used. When configured as four OC-12s or
7 one OC-48, the deserializer 12 converts the serial data stream to
8 a sixteen-bit wide parallel stream. The deserializer 12 includes
9 circuitry to divide the input serial clocks by sixteen. The
10 inputs to the deserializer include a one-bit serial data input, a
11 one-bit 622 MHz clock and a one-bit 155 MHz clock. The outputs
12 include a sixteen-bit parallel data output, a one-bit 38.87 MHz
13 clock and a 9.72 MHz clock. The SONET interfaces are described in
14 more detail in sections 3.2 and 3.3 of Appendix A.

15
16 Parallel data is sent to the SONET framer and transport
17 overhead (TOH) block 14. All incoming signals are framed
18 according to the BELLCORE GR-253 standard which is incorporated
19 herein by reference. The byte boundary and the frame boundary are
20 found by scanning a series of sixteen bit words for the F628
21 pattern. The framer frames on the pattern F6F6F6282828.
22 Independent SONET SPEs within the STS-N frame are demultiplexed by
23 the framer 14. There is a maximum of four independent line
24 interfaces, therefore the framer 14 includes four independent
25 framers. The inputs to the framer include a sixteen-bit parallel

1 data input, and a one-bit clock which will accept 155 MHz, 38.87
2 MHz, or 9.72 MHz. The outputs of the framer include a sixteen-bit
3 parallel data output, a one-bit start of frame (SOF) indication, a
4 six-bit SPE ID used to indicate SONET SPE number. The SPEs are
5 numbered 1 through 48 with respect to the line side port
6 configuration.

7
8 The block 14 also terminates the transport (section and line)
9 overhead for each independent SONET SPE. Since there are a
10 maximum of forty-eight OC-1s on the line side, forty-eight
11 transport overhead blocks are provided unless blocks are time-
12 shared. The inputs to the TOH termination are the same as those
13 discussed above with respect to the framer. The six-bit SPE ID
14 enables data into this block. There is no need for an output data
15 bus as the traffic is routed to this block and to the next block
16 (Ptr Proc 16) on the same data bus. The data path only flows into
17 this block, not through it.

18
19 The pointer processor 16 uses the SONET pointer (H1, H2 and
20 H3 bytes in the TOH) to correctly locate the start of the payload
21 data being carried in the SONET envelope. The SONET pointer
22 identifies the location of byte #1 of the path overhead. The
23 pointer processor 16 is responsible for accommodating pointer
24 justifications that were inserted in order to justify the
25 frequency difference between the payload data and the SONET

1 envelope. Since there are a maximum of forty-eight OC-1s, forty-
2 eight pointer processor blocks are mated to the forty-eight
3 transport overhead termination blocks unless blocks are time-
4 shared. The inputs to the pointer processor 16 are the same as
5 those to the framer and TOH terminator 14. The outputs include a
6 sixteen-bit parallel data output, a one-bit start of SPE indicator
7 which coincides with word 1 of SPE 3, a one-bit SPE valid
8 indicator which gaps out overhead and accommodates pointer
9 movements, and a one-bit POH valid indicator which indicates when
10 a path overhead byte is on the output bus.

11
12 The POH processor 18 processes the nine bytes of Path
13 Overhead in each of the forty-eight SONET SPEs. Since there are a
14 maximum of forty-eight SPEs, forty-eight path overhead processors
15 are provided unless processors are time-shared. The inputs to the
16 path overhead processor 18 include an eight-bit parallel data
17 input, a four-bit SPE ID, the one-bit start of SPE indicator, and
18 the one-bit POH valid indicator. The outputs include a one-bit V1
19 indicator, J1 info, alarms, and path status. Further details
20 about blocks 14, 16, and 18 are provided by the GR-253 standard
21 and documentation accompanying standard SONET mapper/demappers
22 such as those available from Lucent or TranSwitch.

23
24 Once the frame boundaries of the incoming SONET/SDH signals
25 are found and the location of the SPEs has been identified either

1 through pointer processing or through Telecom bus I/F control
2 signals, and the Path Overhead is processed, the payload is
3 extracted from the SPE. The SPEs may be carrying TDM traffic, ATM
4 cells or IP packets. The type of traffic for each SPE is
5 configured through the microprocessor interface 78. Each SPE can
6 carry only one type of traffic. The data from each SPE is routed
7 directly to the correct payload extractor.

8
9 SPEs containing packets and ATM cells are sent to the HDLC
10 framer 20 and the cell delineation block 22, respectively. Each
11 SPE can be configured to carry packet data (packet over SONET).
12 The Port Processor 10 supports packet over SONET for the following
13 SONET (SDH) signals: STS-1 (VC-3), STS-3c (VC-4), STS-12c
14 (VC-4-4c), and STS-48c (VC-4-16c). The datagrams are encapsulated
15 in PPP packets which are framed using the HDLC protocol. The HDLC
16 frames are mapped byte-wise into SONET SPEs and high order SDH
17 VCs. The HDLC framer 20 performs HDLC framing and forwards the
18 PPP packet to a FIFO buffer 24 where it awaits assembly into PDUs.
19 The framer 20 has an input which includes a sixteen-bit parallel
20 data input, a six-bit SPE ID, a one-bit SPE valid indicator, and a
21 one-bit PYLD valid indicator. The output of the framer 20
22 includes a sixteen-bit data bus, a one-bit start of packet
23 indicator, and a one-bit end of packet indicator. Further details
24 about packet extraction from SONET are found in IETF (Internet

1 Engineering Task Force) RFC 1619 (1999) which is incorporated
2 herein by reference.

3
4 The cell delineation block 22 is based on ITU-T G.804, "ATM
5 Cell Mapping into Plesiochronous Digital Hierarch (PDH)", 1998,
6 the complete disclosure of which is hereby incorporated herein by
7 reference. The cell delineation block 22 has inputs which include
8 a sixteen-bit parallel data bus, a six-bit SPE ID, a one-bit SPE
9 valid indicator, and a one-bit POH valid indicator. The outputs
10 include a sixteen-bit parallel data bus and a one-bit start of
11 cell indicator. Cells are placed in a FIFO 24 while awaiting
12 assembly into PDUs. Further details regarding ATM extraction from
13 SONET are found in ITU-T G.804.

14
15 The TDM data is routed to a TDM demultiplexer and low order
16 pointer processor block 26 where the low order VTs and VCs are
17 identified. If a particular SPE is configured for TDM data, then
18 the TDM mapping is described using the host interface 78. Each
19 SPE can carry a combination of VC-11 , VC-12 , VC-2, VC-3 & VC-4.
20 There are seven VT groups in a single STS-1 payload, each VT group
21 has twelve columns. Within one VT Group all of the VTs must be
22 the same. Different VT groups within the same STS-1 SPE can carry
23 different VT types, but within the group it is required that all
24 VTs be of the same type. The VCs and VTs are demultiplexed out of
25 the SONET signal based on the configuration for each of the SPEs.

1 There is no interpretation of the traffic required to locate the
2 containers and tributaries as all of this information is found in
3 the configuration table (not shown) which is configured via the
4 host interface 78. Frames are located inside of the VCs and the
5 VTs through the H4 byte in the path overhead of the SPE. Pointer
6 processing is performed as indicated by the V bytes in the VT
7 superframe. The TDM demultiplexer and low order pointer processor
8 block 26 has inputs which include sixteen bits of parallel data, a
9 six-bits SPE ID, a one-bit start of SPE indicator, a one-bit SPE
10 valid indicator, a one-bit V1 indicator, and one-bit POH valid
11 indicator. The TDM demultiplexer and low order pointer processor
12 block 26 provides the following outputs to the switch mapper 52:
13 sixteen bits of parallel data, a one-bit VT/VC valid indicator, a
14 six-bit SPE ID, and a five-bit VT/VC Number (0-27). The TDM data
15 is placed in reserved slots in the frame as mentioned above and
16 described in more detail below with reference to the switch mapper
17 52. Further details regarding TDM extraction are found in the GR-
18 253 specification

19
20 IP packets and ATM cells from the UTOPIA interface 44 are
21 placed in FIFO 46. Packets and cells from the FIFOs 24 are merged
22 with the packets and cells from the FIFO 46. The descriptor
23 constructor 64 determines whether the data is an ATM cell or an IP
24 packet and generates a corresponding interrupt to trigger the
25 IPF/ATM look-up processor 66 to perform either IP routing look-up

1 or ATM look-up. IP routing look-up is performed by searching for
2 the IP destination address for every packet and the IP source
3 address for packets that need classification. ATM look-up is
4 performed by searching the VPI/VCI fields of the cells. Outputs
5 of the IPF/ATM look-up processor 66 for both IP packets and ATM
6 cells include a seventeen-bit flow index, a five-bit QOS index,
7 and an indicator showing whether the IP packet needs
8 classification. If the IP packet needs classification, the packet
9 is passed to the IP classification processor 68 for
10 classification; otherwise it is passed to the next stage of packet
11 processing, the RED/policing processor 70. IP classification is
12 described in detail in section 6.4 of Appendix A. The
13 RED/Policing processor 70 performs random early detection and
14 weighted random early detection for IP congestion control,
15 performs leaky bucket policing for ATM traffic control, and
16 performs early packet and partial packet discard for controlling
17 ATM traffic which contains packets. The RED/Policing traffic
18 control is described in detail in sections 7.5 et seq. of Appendix
19 A. The presently preferred embodiment of the port processor 10
20 includes a mode register (not shown) which can be placed in a
21 bypass mode to globally turn off the IP/ATM forwarding. In bypass
22 mode, an external device is used for IP/ATM forwarding, and the
23 data descriptors generated by the descriptor constructor 64 are
24 routed directly to an output FIFO (not shown).

25

1 All of the data stored in the FIFOs 24 and 46 is in fifty-
2 two-byte "chunks". If an IP packet is longer than fifty-two-
3 bytes, it is segmented into multiple fifty-two-byte chunks. The
4 input data descriptor for each chunk includes indications of
5 whether the chunk is an ATM cell or a packet, whether it is the
6 start of a packet or the end of a packet, packet length, and the
7 source and destination port numbers. After processing by the
8 IPF/ATM lookup processor 66 and the IP classification processor
9 68, an output data descriptor is written to a FIFO (not shown)
10 which is read by the RED/Policing processor 70.

11
12 Cells and packets which survive RED/policing are read by the
13 receive data link manager 72 which creates the PDUs described
14 above with reference to Figure 3a. The receive data link manager
15 is described in detail in section 8 of Appendix A. According to
16 the presently preferred embodiment, processed cells and packets
17 are stored in an external FIFO which is read whenever it is not
18 empty.

19
20 As shown in Figure 1, the switch mapper 52 receives TDM
21 traffic from the TDM demultiplexer and low order pointer processor
22 26 as well as PDUs from the data link manager 72. As mentioned
23 above, the switch mapper also receives request elements. The
24 request elements are formed by the arbiter 56 as described in more
25 detail below. It is the function of the switch mapper (also

1 referred to as the data mapper in Appendix A) to arrange TDM data,
2 PDUs, and request elements in the frame described above with
3 reference to Figures 3 and 3a-c.

4
5 The switch mapper 52 includes a state machine (not shown)
6 which is associated with the ATM/IP PDUs. The data link manager
7 72 writes the PDU's using a sixty-four-bit interface to the
8 external FIFO (not shown). The data is transmitted from the
9 external FIFO to the switch mapper 52 in thirty-two-bit slots with
10 four bits of parity. The state machine associated with the
11 external PDU FIFO monitors the status of the FIFO and maintains
12 data integrity.

13
14 In section 9 of Appendix A, the data link manager 72, arbiter
15 block 56, switch mapper 52, and weighted round robin scheduler 80,
16 together with memory and other support circuits (not shown in
17 Figure 1) are referred to collectively as the "receive switch
18 controller". As described in detail above, each incoming ATM cell
19 and packet is processed by performing a lookup based on the ATM
20 VPI/VCI or on the IP source and destination. This lookup first
21 verifies that the connection is active, and if active, it returns
22 a seventeen-bit index. For ATM cells, the index points to a set
23 of per VC parameters and to routing information. For packets, the
24 index points to a set of queuing parameters and to routing
25 information. The seventeen-bit index supports a maximum of 128K

1 simultaneous IP and ATM flows through the port processor. The ATM
2 cells are encapsulated in a cell container and stored in one of
3 128K queues in external memory. These 128K queues are managed by
4 the data link manager 72. As mentioned above, the IP packets are
5 fragmented into fifty-two-byte blocks and each of these blocks is
6 encapsulated in a cell container (PDU). These cell containers are
7 also stored in one of the 128K queues in external memory by the
8 data link manager. The 128K IP/ATM flows are aggregated into one
9 of thirty-two QOS queues for scheduling through the switch. The
10 data link manager 72 also aggregates all the control headers
11 required for transmission of cells through the switch into the QOS
12 queues and inserts these routing tags into one of thirty-one QOS
13 routing tag FIFOs. One of the queues, is reserved for high
14 priority traffic. Any cells arriving in the high priority queue
15 will interrupt the scheduler 80 and will be scheduled to leave the
16 high priority queue immediately.

17
18 The scheduler 80 is responsible for scheduling cell
19 containers through the switch. The scheduling algorithm used is
20 weighted round robin which operates on the QOS queues. Once cells
21 have been scheduled from these queues, the control headers from
22 these queues are forwarded to the arbiter 56 and are stored in a
23 request control table (not shown). The request arbiter 56 forms
24 request elements from the control headers and forwards these
25 requests to the switch data mapper 52 for transmission through the

1 switch. The grants received in response to these requests are
2 deserialized by block 58, deframed and transferred back to the
3 arbiter block 56 by the grant block 62. For granted requests, the
4 cell containers are dequeued from external memory by the data link
5 manager 72 and transferred to the switch mapper 52 for
6 transmission through the switch.

7
8 As mentioned above, the port processor 10 supports redundancy
9 in order to improve reliability. Two redundancy schemes are
10 supported. In the first redundancy scheme, the switch controller
11 supports redundant routing tags and transparent route switch-over.
12 In the second redundancy scheme, the port processor supports
13 redundant data channels in both input and output directions.
14 The redundant data channels connect to two separate switch
15 fabrics. In the Appendices they are referred to as the A and B
16 data channels. Each control header contains two routing tags, and
17 each routing tag has a corresponding AB channel tag. This
18 provides for two routes through the switch for data transmission.
19 If both routing tags have the same channel tag, this allows for
20 two alternate paths through the same switch fabric. If both
21 routing tags have different channel tags, this allows for a
22 redundant switch fabric and any route failing in one switch fabric
23 will cause a switch-over to use the redundant switch fabric. An
24 AB channel tag is used to indicate whether the data is to be
25 routed using the A data channel or the B data channel. If, after

1 a programmable number of consecutive tries, no grant is received
2 in response to request elements using the A channel routing tag, a
3 bit is set to switch over to the B channel routing tag. Further
4 details of the redundancy feature are provided in sections 10.2.3
5 and 9.2.3 of Appendix A.

6
7 As mentioned above, the arbiter 56 is responsible for sending
8 requests to the switch mapper 52 and processing the grants that
9 arrive from the grant demapper 62. The arbiter dequeues requests
10 from a routing tag FIFO, copies this information into a request
11 control table, writes the FLOWID into FLOWID RAM, resets a request
12 trial counter that counts the number of times a request has been
13 tried, and resets the grant bit. Each request message has a
14 unique request ID which is returned in the grant message. The
15 request ID is the index in the arbiter request control table into
16 which the routing tag is copied. The routing tag along with the
17 request ID is forwarded to a routing tag formatter block which
18 formats the routing tag into a request message and inserts the
19 request into a request FIFO in the switch mapper 52.

20
21 The grant demapper in the grant block 62 stores the request
22 ID and the grant in a FIFO called the grant_reqid FIFO. In the
23 arbiter block 56, the request IDs are dequeued from A and B
24 grant_reqid FIFOs alternatively depending on whether the
25 switchover bit is set. The request IDs dequeued from the FIFO are

1 used to set a grant bit in the grant register at the bit position
2 indicated by the request ID, to index the FLOWID RAM, and read the
3 FLOWID associated with the request ID. This FLOWID is written
4 into a deq-flowid FIFO for the appropriate channel, i.e. if the
5 request ID is dequeued from the A reqid_fifo, the FLOWID is
6 written into the A deqflowid_fifo. The data link manager 72
7 monitors the deqflowid_fifo and uses the FLOWID to dequeue data
8 PDUs from external memory and send them to the switch mapper 52
9 for transmission in the next row time.
10

11 An end_of_grants signal is asserted by the grant demapper 62,
12 when no more grants can be received at the grant demapper. In
13 most switch implementations the end_of_grants signal is rarely, if
14 ever, asserted. It is only in switches having many stages that
15 the end_of_grants signal is more likely to be asserted. Once the
16 end_of_grant signal has been received the arbiter 56 begins the
17 process of updating the request control table. If a grant has not
18 been returned for a routing tag stored in the request control
19 table, the request trial counter is incremented and a new request
20 is generated using the routing tag. If a routing tag in the
21 request control table has been sent as a RE a (programmed) maximum
22 number of times, the most significant fifteen bits of the FLOWID
23 are used to index into the redundancy control table and update the
24 bit to indicate failure of the current path and to select the

1 alternate routing path. Further details regarding the arbiter
2 block 56 are provided at section 9.2.4 of Appendix A.

3
4 As described above, the TDM data, ATM/IP PDU's and the
5 request messages are combined into a single data stream for
6 transmission through the switch fabric. This combination is
7 performed by the switch mapper 52 on the receive side of the port
8 processor. On the transmit side of the port processor, a switch
9 demapper 60 separates TDM data from ATM/IP PDUs. According to the
10 presently preferred embodiment, the demapper 60 is provided with
11 external memory for a PDU FIFO. For ATM/IP data, the demapper
12 writes PDUs to the FIFO and interrupts the data link manager 74.
13 The data link manager 74 reads the header information from the PDU
14 FIFO, and extracts the FLOWID. Based on the FLOWID, the datalink
15 manager 74 retrieves a Linked List/Shaping/Scheduling data
16 structure from external memory. The data link manager 74 writes
17 the linked list pointers to the PDU FIFO, then initiates a DMA
18 transfer to move the PDU to external memory. The data link
19 manager updates the head, tail, and count fields in the Linked
20 List/Shaping/Scheduling data structure and passes the data
21 structure to the Shaping/Scheduling processor 76 through a
22 Shaping/Scheduling FIFO. The Shaping/Scheduling processor 76
23 performs the Shaping and Scheduling functions and updates the
24 Linked List/Shaping/Scheduling datastructure.

25

1 The data flow from external memory to the SONET/UTOPIA Data
2 FIFOs 30 and 48 is as follows. The data link manager 74 polls the
3 PDU FIFO and SONET/UTOPIA FIFO status flags. If the PDU FIFO is
4 not empty and the SONET/UTOPIA FIFO is not full for a particular
5 output port, the data link manager 74 retrieves the Link
6 List/Shaping/Scheduling data structure for the Flow ID read from
7 the PDU FIFO. (Note that for an IP packet flow, the data link
8 manager will continue to retrieve PDUs from the Linked List until
9 a PDU with an End of Packet indicator is found.) The data link
10 manager then initiates a DMA transfer from external memory to the
11 SONET/UTOPIA FIFOs 30, 48. The data link manager 74 then updates
12 the Link List/Shaping/Scheduling data structure and writes it back
13 to external memory.

14
15 On the transmit side of the port processor 10, the grant
16 framer, deframer, serializer and deserializer in the grant block
17 62, the switch demapper 60, the transmit datalink manager 74, and
18 the transmit scheduler and shaper 76 are referred to collectively
19 as the transmit (TX) switch controller in Appendix A. The TX
20 switch controller is responsible for either accepting or rejecting
21 requests that come into the port processor for output
22 transmission. To do this, the TX switch controller checks if the
23 queue identified by the output port number of the request can
24 accept a cell container. These one hundred twenty-eight queues
25 are managed by the TX data link manager 74. According to the

1 presently preferred embodiment, these queues are stored in
2 external memory. The scheduling of these cell containers is
3 performed by the TX scheduler 76. If the queue can accept the
4 cell container, the request is turned into a grant and inserted
5 into a grant_fifo. The grant-framer and serializer 62 reads this
6 information and creates an grant message for transmission through
7 the grant path.

8
9 The TX switch controller monitors the status of the data
10 queues for each of the one hundred twenty-eight output ports using
11 the following three rules. If the full_status bit for the
12 requested output port is set, there is no buffer space in the
13 queue for any data PDUs destined for that output port and all
14 requests to that output port are denied. If the full_status bit
15 is not set and the nearly_full_status bit is set, there is some
16 space in the queue for data PDUs destined for that output port;
17 however this space may be reserved for higher priority traffic.
18 In this instance the QOS number is checked against a threshold
19 (programmed) QOS number and if the QOS number is less than the
20 threshold, the request will be accepted. If the nearly
21 full_status bit is not set, all incoming requests are granted. If
22 a request is accepted, the corresponding output port counter is
23 incremented. This reserves space in the data buffer (30 or 48)
24 for the arrival of the data PDU at that output port. The transmit
25 data link manager 74 constantly monitors the one hundred twenty-

1 eight output port counters and sets/resets the one hundred twenty-
2 eight full and nearly full status bits.

3
4 The port processor 10 creates complete outgoing SONET
5 signals. All of the transport and path overhead functions are
6 supported. The SONET interfaces can run in source timing mode or
7 loop timing mode.

8
9 The high order pointer is adjusted by the high order pointer
10 generator 38 through positive and negative pointer justifications
11 to accommodate timing differences in the clocks used to generate
12 the SONET frames and the clock used to generate the SONET SPEs.
13 At initialization, SPE FIFOs are allowed to fill to halfway before
14 data is taken out. The variations around the center point are
15 monitored to determine if the rate of the SONET envelope is
16 greater than or less than the rate of the SPE. If the rate of the
17 SONET envelope is greater than the rate of the SPE, then the SPE
18 FIFO will gradually approach a more empty state. In this case,
19 positive pointer movements will be issued in order to give the SPE
20 an opportunity to send additional data. If the rate of the SONET
21 envelope is less than the rate of the SPE, then the SPE FIFO will
22 gradually approach a more full state. In this case, negative
23 pointer movements will be issued in order to give the SPE an
24 opportunity to output an extra byte of data from the FIFO. The

1 SONET framer and TOH generator 40 generate transport overhead
2 according to the BELLCORE GR-253 standard.

3
4 The outgoing SONET frames are generated from either the
5 timing recovered from the receive side SONET interface or from the
6 source timing of the Port Processor. Each signal is configured
7 separately and they can be configured differently. The frame
8 orientation of the outgoing SONET frames is arbitrary. Each of
9 the four signals can be running off different timing so there is
10 no need to try to synchronize them together as they will
11 constantly drift apart. There is no need to frame align the Tx
12 ports to the Rx ports as this would result in realigning the Tx
13 port after every realignment of the Rx port.

14
15 For OC-3 and OC-12 the 16-bit wide internal bus is serialized
16 to 155 Mbps or 622 Mbps by the serializer 42. For OC-48
17 applications, the entire sixteen bit bus is output under the
18 control of an external serializer (not shown).

19
20 There is a potential for forty-eight different SPEs being
21 generated for the outgoing SONET interfaces. All of these SPEs
22 are generated from a single timing reference. This allows all of
23 the SPE generators to be shared among all of the SONET and Telecom
24 bus interfaces without multiplexing between the different clocks

1 of the different SONET timing domains. The SPE consists of the
2 Path level overhead and the payload data. The payload data can be
3 TDM, ATM or packet. All of these traffic types are mapped into
4 single SPEs or concatenated SPEs as required by their respective
5 standards. As the SPEs are generated, they are deposited into SPE
6 FIFOs. For each SPE there is a sixty-four-byte FIFO and these
7 individual SPE FIFOs are concatenated through SPE concatenation
8 configuration registers. As described above, the fill status of
9 the SPE FIFOs is used to determine the correct time to perform a
10 positive or negative pointer justification.

11
12 TDM, ATM and packet data are all mapped into SONET SPEs as
13 specified by their respective standards. The type of data carried
14 in each of the potential forty-eight SPEs is configured through
15 the external host processor. Based on this configuration, each
16 SPE generator is allocated the correct type of mapper. All of
17 this configuration is performed at initialization and can only be
18 changed when the particular SPE is first disabled. Once the
19 configuration is complete, there is an isolated set of
20 functional blocks allocated to each SPE. This set of functional
21 blocks includes one of each of the following: payload mapper,
22 payload FIFO, POH generator, SPE FIFO and SPE generator.
23 Each of the ATM and packet payload mappers has a payload FIFO into
24 which it writes payload data for a particular SPE. For TDM

1 traffic, each potential Virtual Container is allocated its own
2 FIFO.

3
4 Returning now to Figure 2, in each "datapath and link
5 bandwidth arbitration module" 102, the data stream deserializer
6 126 synchronizes to the incoming serial data stream and then
7 reassembles the row stream which is transported using two physical
8 unilink channels. It also provides FIFO buffering on each of the
9 two incoming serial streams so that the streams may be "deskewed"
10 prior to row reassembly. It recovers the thirty-six-bit slot data
11 from the row stream in a third FIFO which is used for deskewing
12 the twelve input links. This deskewing allows all the input links
13 to forward slot N to the switching core simultaneously. The link
14 deskewing is controlled by the link synchronization and timing
15 control module 150. The deserializer 126 also continuously
16 monitors the delta between where slot 0 of the incoming row is
17 versus the internal row boundary signal within the switch element.
18 The difference is reported to the Link RISC Processor 156 and is
19 used (in the first stage of a switch) as part of the ranging
20 process to synchronize the port processor connected to the input
21 link.

22
23 The data stream demapper 128 is responsible for extracting
24 the data from the incoming serial data links. It demaps the input
25 link slots based on the input slot number and determines whether

1 the traffic is TDM, PDU, or a request element (RE). For TDM
2 traffic, the demapper determines the destination link and row
3 buffer 132 memory address. This information is stored in a
4 demapper RAM (not shown) which is configured by software when TDM
5 connections are added or torn down. For PDU traffic, the demapper
6 128 assembles all sixteen slots which make up the PDU into a
7 single 64-byte PDU word, then forwards this entire PDU word to the
8 row buffer mapper logic 130. The PDUs are assembled prior to
9 forwarding them to the row buffer 132 so that the row buffer
10 mapper 130 can write the entire PDU to the row buffer 132 in a
11 single clock cycle. This provides the maximum possible write-side
12 memory bandwidth to the row buffer 132. It is a significant
13 feature of the switch element that twelve entire PDUs are written
14 to a single row buffer in six link slot times (twelve core clock
15 cycles). For request elements, the demapper 128 assembles the
16 three-slot block of REs into two forty-eight-bit REs and forwards
17 them to the request parser module 152. A detailed description of
18 the data stream demapper 128 is provided in Sections 4.3.1 et seq.
19 of Appendix B.

20
21 The row buffer mapper 130 is responsible for mapping traffic
22 which is received from the data stream demapper 128 into the row
23 buffer 132. The mapper 130 provides FIFO buffers for the TDM
24 traffic as it is received from the data stream demapper 128, then

1 writes it to the row buffer 132. The row buffer memory address is
2 actually preconfigured in the demapper RAM (not shown) within
3 the data stream demapper module 128. That module forwards the
4 address to the row buffer mapper 130 along with the TDM slot data.
5 The mapper 130 also writes PDU traffic from the data stream
6 demapper 128 to the row buffer 132 and computes the address within
7 the row buffer 132 where each PDU will be written. PDUs are
8 written into the row buffers starting at address 0 and then every
9 sixteen-slot address boundary thereafter, up to the maximum
10 configured number of PDU addresses for the row buffer 132. A
11 detailed description of the row buffer mapper 130 is provided in
12 Section 4.3.1.4 of Appendix B.

13
14 The row buffer 132 contains the row buffer memory elements.
15 According to the presently preferred embodiment, it provides
16 double buffered row storage which allows one row buffer to be
17 written during row N while the row data which was written during
18 row N-1 is being read out by the data stream mapper 136. Each row
19 buffer is capable of storing 1536 slots of data. This allows the
20 row buffer to store ninety-six PDUs or 1536 TDM slots or a
21 combination of the two traffic types. Request elements and link
22 overhead slots are NOT sent to the row buffer 132. Therefore the
23 row buffer does not need to be sized to accommodate the entire
24 1700 input link slots. According to the presently preferred
25 embodiment, the row buffer write port is $16 \times 36 = 576$ bits wide and

1 it supports writing of only one thirty-six-bit slot (TDM data) or
2 writing of an entire 576-bit word (PDU data) in a single clock
3 cycle. A detailed description of the row buffer 132 is provided
4 in Section 4.3.1.4 of Appendix B.

5
6 Request arbitration utilizes two components: a centralized
7 request parser module 152 and a request arbitration module 134 for
8 each of the output links. Request elements are extracted from the
9 input slot stream by the data stream demapper 128 and are
10 forwarded to the request parser 152. The request parser 152
11 forwards the forty-eight-bit request elements to the appropriate
12 request arbitration module 134 via two request buses (part of the
13 input link bus 120). Each request bus may contain a new request
14 element each core clock cycle. This timing allows the request
15 arbitration logic to process thirteen request sources in less than
16 eight core clock cycles. The thirteen request sources are the
17 twelve input data streams and the internal multicast and in-band
18 control messaging module 156. The request arbitration module 134
19 monitors the two request element buses and reads in all request
20 elements which are targeted for output links the request
21 arbitration module is implementing. According to the presently
22 preferred embodiment, the request arbitration module 134 provides
23 buffering for up to twenty-four request elements. When a new
24 request element is received, it is stored in a free RE buffer (not
25 shown). If there are not any free RE buffers, then the lowest

1 priority RE which is already stored in a buffer is replaced with
2 the new RE if the new RE is a higher priority. If the new RE is
3 equal to or lower in priority than all REs currently stored in the
4 RE buffers then the new RE is discarded. On the output side, when
5 the data stream mapper module 138 is ready to receive the next RE,
6 the request arbitration module 134 forwards the highest priority
7 RE which is stored in the RE buffers to the data stream mapper
8 module 136. If the RE buffers are empty, then an "Idle" RE is
9 forwarded. A detailed description of the request arbitration
10 module 134 is provided in Section 7 of Appendix B.

11
12 The data stream mapper 136 is responsible for inserting data
13 and request elements into the outgoing serial data links. This
14 includes mapping of the output link slots based on the output slot
15 number to determine if the traffic is TDM, PDU, request element,
16 or test traffic. The determination is based on the contents
17 of the mapper RAM (not shown). For TDM traffic, the row buffer
18 memory address is determined from the mapper RAM which is
19 configured by software as TDM connections are added or torn down.
20 For PDU traffic, the data stream mapper 136 reads one slot at a
21 time from the row buffer 132. The row buffer memory address is
22 stored in the mapper RAM by software. If the target PDU is not
23 valid (i.e., a PDU was not written to that row buffer location
24 during the previous row time), then the mapper 136 transmits an
25 idle pattern in order to ensure that a data PDU is not duplicated

1 within the switch. For request elements, the mapper assembles the
2 three-slot block of REs from two forty-eight-bit REs. The REs are
3 read from the request arbitration module 134. For test patterns,
4 the mapper 136 inserts the appropriate test pattern from the
5 output link bus 122. These test patterns are created by either
6 the test pattern generator 162 or test interface bus 164 modules.

7
8 The data stream mapper supports slot multicasting at the
9 output stage. For example, the data stream mapper for any output
10 link is able to copy whatever any other output link is sending out
11 on the current slot time. This copying is controlled via the
12 mapper RAM and allows the mapper to copy the output data from
13 another output link on a slot-by-slot basis. A detailed
14 description of the data stream mapper 136 is provided in Section 4
15 of Appendix B.

16
17 The data stream serializer 138 creates the output link serial
18 stream. Data slots are received via the data stream mapper module
19 136 and the link overhead is generated internally by the data
20 stream serializer 138. The serializer 138 also splits the row
21 data stream into two streams for transmission on the two paths
22 110, 114. A detailed description of this module is provided in
23 Section 11 of Appendix B.

1 The grant stream deserializer 140 in each module 102 works in
2 much the same manner as the data stream deserializer 126. The
3 primary difference is that the grant data only utilizes a single
4 path, thus eliminating the need for deskewing and deinterleaving
5 to recover a single input serial stream. Since this serial link
6 is only one half the data stream rate of the forward link, there
7 are 850 slots per row time. A single FIFO (not shown) is used to
8 allow for deskewing of the input serial grant streams for all 12
9 links. A detailed description of the grant stream deserializer
10 140 is provided in Section 11 of Appendix B.

11
12 The grant stream demapper 142 is responsible for extracting
13 the data from the incoming serial grant links. This includes
14 demapping of the received grant link slots based on the input slot
15 number to determine if the traffic is a grant element or another
16 kind of traffic. The determination is based on the contents of
17 the grant demapper RAM (not shown). According to the presently
18 preferred embodiment, traffic other than grant elements is not yet
19 defined. For grant elements, the grant stream demapper 142
20 assembles the three-slot block of GEs into two forty-eight-bit GEs
21 and forwards them to the single grant parser module 154. A
22 detailed description of the grant stream demapper 142 is provided
23 in Section 7.2.3.2 of Appendix B.

24

1 The grant arbitration module 144 operates in an identical
2 manner to the request arbitration logic 134. In the presently
3 preferred embodiment, this module is identical to the request
4 arbitration module. The only difference is that it processes
5 grant elements in the reverse path instead of request elements in
6 the forward path. It will be recalled that grant elements are, in
7 fact, the request elements which have been returned.

8
9 The grant stream mapper 146 is responsible for inserting data
10 into the outgoing serial grant links. It maps the output grant
11 slots based on the output slot number to determine if the
12 traffic is a grant element or test traffic. The determination is
13 based on the contents of the grant mapper RAM (not shown). For
14 grant elements, it assembles the three-slot block of GEs from two
15 forty-eight-bit GEs. The GEs are read from the grant arbitration
16 module 144. For test patterns, it inserts the appropriate test
17 pattern from the output link bus 122. These test patterns are
18 created by either the test pattern generator 162 or the test
19 interface bus 164 modules. A detailed description of the grant
20 stream mapper 146 is provided in Section 7.2.3.2.

21
22 The grant stream serializer 148 works in much the same manner
23 as the data stream serializer 138. The primary difference is that
24 the grant data only utilizes a single path, thus eliminating the
25 need for interleaving the transmit serial stream across multiple

1 output serial streams. Since this serial link is only one half
2 the forward data stream rate, there are only 850 slots per row
3 time. A detailed description of the grant stream serializer 148
4 is provided in Section 11 of Appendix B.

5
6 The modules described above (except for the request parser
7 and the grant parser) are instantiated for each link module 102 of
8 which there are twelve for each switch element 100. The following
9 modules are instantiated only once for each switch element.

10
11 The link synchronization & timing control 150 provides the
12 global synchronization and timing signals used in the switch
13 element. It generates transmission control signals so that all
14 serial outputs start sending row data synchronized to the RSYNC
15 (row synchronization) input reference. It also controls the
16 deskewing FIFOs in the data stream deserializers so that all
17 twelve input links will drive the data for slot N at the same time
18 onto the input link bus 120. This same deskewing mechanism is
19 implemented on the grant stream deserializers. A detailed
20 description of the link synchronization and timing control 150 is
21 provided in Section 10 of Appendix B.

22
23 The request parser 152 receives inputs from all thirteen
24 request element sources and forwards the REs to the appropriate
25 request arbitration modules via the two request element buses. A

1 detailed description of the request parser 152 is provided in
2 Section 7.2.1.1 of Appendix B.

3
4 The grant parser 154 physically operates in an identical
5 manner to and is identical to the request parser 152. The only
6 difference is that it processes grant elements in the reverse path
7 instead of request elements in the forward path. As mentioned
8 above, the grant elements contain the same information as the
9 request elements, i.e. the link address through the switch from
10 one port processor to another.

11
12 The link RISC processor 156 controls the ranging
13 synchronization on the input links with the source port processors
14 in the first stage of the switch fabric. It also controls the
15 ranging synchronization on the output link grant stream input with
16 the source port processors in the last stage of the switch fabric.
17 It also handles the Req/Grant processing needed to transmit
18 multicast messages and controls the reception and transmission of
19 the in-band communications PDUs. All in-band communications PDUs
20 are forwarded to the Configuration RISC Processor 158 which
21 interprets the messages. The link RISC processor 156 only handles
22 the Req/Grant processing needed to transmit multicast and in-band
23 communications messages.

24

1 The configuration RISC controller 158 processes configuration
2 and status messages received from an external controller module
3 (not shown) and in-band communication messages as described above.
4 The system control module 160 handles all the reset inputs and
5 resets the appropriate internal modules. The configuration RISC
6 controller 158 and the system control module 160 are preferably
7 implemented with an Xtensa™ processor from Tensilica, Inc., Santa
8 Clara, CA.

9
10 The test pattern generator and analyzer 162 is used for the
11 generation of various test patterns which can be sent out on
12 any slot on the data stream or grant stream outputs. It is also
13 capable of monitoring input slots from either the received data
14 stream or grant stream.

15
16 The test interface bus multiplexer 164 allows for sourcing
17 transmit data from the external I/O pins and forwarding data to
18 the I/O pins. This is used for testing the switch element when a
19 port processor is not available.

20
21 The unilink PLL 166 is used to create the IF clock needed by
22 the unilink macros. Within each unilink macro another PLL
23 multiplies the IF clock up to the serial clock rate. The core PLL
24 168 is used to create the clock used by the switch element core
25 logic. In the presently preferred embodiment, the core clock is

1 approximately 250 MHz. A detailed description of both PLLs is
2 provided in Section 9 of Appendix B.

3
4 The JTAG interface 170 is used for two purposes: (1) boundary
5 scan testing of the switch element at the ASIC fab and (2) Debug
6 interface for the Configuration RISC Processor.

7
8 As shown in Figure 2, there are three datapath buses (the
9 input link bus 120, the output link bus 122, and the grant bus
10 124) which are used to move switched traffic from the input links
11 to the output links. These buses are also used to carry traffic
12 which is sourced or terminated internally within the switch
13 element. The significant datapaths of the input link bus are
14 summarized in Table 2 below. The significant datapaths of the
15 output link bus are summarized in Table 3 below. The significant
16 datapaths of the grant bus are summarized in Table 4 below.

Name	Qty	Width	Description	Source
islot_num	1	11	Current input slot number for traffic from the Data Stream Deserializers	Link Sync & Timing Ctrl
ilink_req_0 thru ilink_req_11	12	48	Request elements received on the input link	Data Stream Demapper module for each input link
lcl_req_0	1	48	Request elements generated locally	Link RISC Controller
req_a, req_b	2	48	Parsed request elements	Request Parser
ilink_tdm_data_0 thru ilink_req_11	12	47	TDM data, 36-bit data + 11 bit destination row buffer address	Data Stream Demapper module for each input link.
ilink_tdm_dlink_0 thru ilink_tdm_dlink_11	12	4	Destination output link (i.e., row buffer) identifier	Data Stream Demapper module for each input link
ilink_pdu_0 thru ilink_pdu_11	12	512	Complete 64-byte PDU which has been assembled from the incoming slots	Data Stream Demapper module for each input link
ilink_pdu_flag_0 thru ilink_pdu_flag_11	12	13	Each flag is asserted for each destination which the current PDU is addressed. Total destinations = 12 output links plus the internal MC & In-band Comm Controller	Data Stream Demapper module for each input link
lcl_pdu	1	64	Bus used to transport locally generated PDUs to the Data Stream Demappers	Link RISC Controller

Table 2

Name	Qty	Width	Description	Source
oslot_num	1	11	Current output slot number for traffic destined for the output links.	Link Sync & Timing Ctrl
rbuf_dout_0 thru rbuf_dout_11	12	48	Slot data output from the row buffer.	Row Buffer for each output link.
rbuf_rd_addr	12	12	Row buffer read address.	Data Stream Mapper for each output link.
test_src1, test_src2, test_src3	3	36	Test traffic sources.	Test Pattern Generator, Test Interface Bus
idle_ptrn	1	36	Idle pattern which is transmitted when no valid PDU data is available.	Data Stream Demapper module for each input link.
olink_req_0 thru olink_req_11	12	48	Request elements for each output link.	Req Arbitration modules.
omap_data_0 thru omap_data_11	12	36	Link output after the mapping multiplexers. All 12 outputs are fed back into each of the Data Stream Mappers so that TDM multicasting can be done.	Data Stream Mapper for each output link

Table 3

Name	Qty	Width	Description	Source
olink_gntslot_num	1	10	Current input slot number for traffic from the Grant Stream Deserializers.	Link Sync & Timing Ctrl
olink_gnt_0 thru olink_gnt_11	12	48	Grant elements received on the grant receiver which is associated with the output link.	Grant Stream Demapper.
olink_gntslot_0 thru olink_gndslot_11	12	36	Demapped slots from the received grant stream. These are slots which are not carrying grant elements.	Grant Stream Demapper.
gnt_a, gnt_b	2	48	Parsed grant elements	Grant Parser

Table 4

According to the presently preferred embodiment, each switch element includes a multicast controller and a separate multicast PDU buffer. Multicast request elements flow through the switch in the same manner as standard unicast request elements. At the point where the message needs to be multicast, the hop-by-hop field's bit code for that switch stage indicates that the request is multicast. The request is forwarded to the multicast controller. On the grant path, the multicast controller sources a grant if there is room for the data in the multicast recirculating buffers. Once the data has been transmitted to the multicast buffer, the multicast controller examines the data header and determines which output links it needs to be sent out on. At this point, the multicast controller sources a number of request messages which are handled in the same manner as unicast requests.

1

2 There have been described and illustrated herein several
3 embodiments of a network switch which supports tdm, atm, and ip
4 traffic. While particular embodiments of the invention have been
5 described, it is not intended that the invention be limited
6 thereto, as it is intended that the invention be as broad in scope
7 as the art will allow and that the specification be read likewise.

8 For example, while the invention has been disclosed as sending
9 requests in-band and sending grants out-of-band, it is within the
10 scope of the claimed invention to send requests out-of-band if
11 bandwidth is constrained. It will therefore be appreciated by
12 those skilled in the art that yet other modifications could be
13 made to the provided invention without deviating from its spirit
14 and scope as so claimed.

15